



# Querying Attributed DL-Lite Ontologies Using Provenance Semirings

Camille Bourgaux, Ana Ozaki

## ► To cite this version:

Camille Bourgaux, Ana Ozaki. Querying Attributed DL-Lite Ontologies Using Provenance Semirings. Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19), Jan 2019, Honolulu, United States. hal-02109645

**HAL Id: hal-02109645**

**<https://inria.hal.science/hal-02109645>**

Submitted on 25 Apr 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Querying Attributed DL-Lite Ontologies Using Provenance Semirings

**Camille Bourgaux**

Télécom ParisTech & DI ENS, CNRS, ENS,  
PSL University & Inria, France

**Ana Ozaki**

KRDB Research Centre,  
Free University of Bozen-Bolzano, Italy

## Abstract

Attributed description logic is a recently proposed formalism, targeted for graph-based representation formats, which enriches description logic concepts and roles with finite sets of attribute-value pairs, called annotations. One of the most important uses of annotations is to record provenance information. In this work, we first investigate the complexity of satisfiability and query answering for attributed DL-Lite<sub>R</sub> ontologies. We then propose a new semantics, based on provenance semirings, for integrating provenance information with query answering. Finally, we establish complexity results for satisfiability and query answering under this semantics.

## Introduction

Description logic (DL) (Baader et al. 2007) ontologies allow to express complex relationships between concepts and roles, but they are ill-equipped to represent and reason about multiple and heterogeneous types of meta-knowledge, such as the temporal validity of a fact, or its source. For instance, the YAGO ontology (Hoffart et al. 2013) attaches provenance metadata to its facts (e.g., source and confidence of the extraction) as well as temporal and geospatial information. Many practical applications therefore use *knowledge graphs*, which consist, like DL assertions, of directed labelled graphs but that also allow, unlike DLs, to add annotations to vertices and edges. Property Graph, the data model used in many graph databases (Rodriguez and Neubauer 2010), and Wikidata, the knowledge graph used by Wikipedia (Vrandečić and Krötzsch 2014), are prominent examples of such labelled graphs. To bridge the gap between DL and knowledge graphs, *attributed description logics* (Krötzsch et al. 2017; Krötzsch et al. 2018) have been recently introduced. They enrich DL concepts and roles with finite sets of attribute-value pairs, called *annotations*, and allow to express constraints on these annotations in the ontology inclusions. For example, the attributed DL assertion  $\text{spouse}(\text{taylor}, \text{burton})@[\text{start} : 1975, \text{end} : 1976]$  states that Liz Taylor was married to Richard Burton from 1975 to 1976, and the following role inclusion expresses that spouse is a symmetric relation, where the inverse statement has the same start and end dates:

$$\text{spouse}@X \sqsubseteq \text{spouse}^-@[\text{start} : X.\text{start}, \text{end} : X.\text{end}].$$

While the work by Krötzsch et al. studied the complexity of the satisfiability problem for several attributed DL languages,

our focus in this paper is on *query answering* in attributed DL. The problem of querying DL ontologies using database-style queries (in particular, conjunctive queries) is an important research topic for which tractable DL languages have been tailored (Bienvenu and Ortiz 2015). We consider here the DL-Lite<sub>R</sub> dialect of the *DL-Lite family* (Calvanese et al. 2007), which underlies the OWL 2 QL profile (Motik et al. 2009), and investigate attributed DL-Lite<sub>R</sub>.

One of the main motivations of attributed DLs is to integrate annotations carrying provenance information, which are very frequent in knowledge graphs<sup>1</sup>. Recording and tracking provenance information is an important topic in database theory, where *provenance semirings* (Green, Karvounarakis, and Tannen 2007) were introduced as an abstract tool to relate the result of a query with information about the original sources of the data and the ways in which the query was obtained. Such information comes in the form of a *provenance polynomial*. It has been useful for many applications, such as query answer explanation or querying of probabilistic databases (Senellart 2017; Cheney, Chiticariu, and Tan 2009; Suciu et al. 2011). Bienvenu, Deutch, and Suchanek (2012) argued that provenance would be useful for Web data, e.g., to establish the authorship or determine the trust in a given piece of data, or to help to guarantee the privacy of information. Provenance has also been investigated for non-relational databases and Semantic Web (see Conclusion for discussion of related work). In this work, we propose a new semantics for the attributed DL annotations, based on provenance semirings, so that queries can be annotated with provenance polynomials. To the best of our knowledge, this is the first work where provenance polynomials are embedded into both the syntax and the semantics of the query.

The first section introduces attributed DL-Lite<sub>R</sub>, following the formalism given by Krötzsch et al. (2017; 2018). We then define attributed conjunctive queries and study the complexity of satisfiability and query answering in attributed DL-Lite<sub>R</sub>. We next present our new semantics for the annotations to model provenance and analyse the complexity of satisfiability and query answering with this new model, considering queries that can be annotated with provenance polynomials. In particular, we show that satisfiability and query answering

<sup>1</sup>E.g., in Wikidata *reference* (provenance) is one the most frequent types of annotations <https://www.wikidata.org>.

in attributed DL-Lite<sub>R</sub> are PSPACE-complete problems. For the semirings-based semantics and queries annotated with provenance polynomials, we establish that although satisfiability is EXPTIME-hard in the general case, the new semantics does not increase the complexity of query answering if the ontology contains only assertions and a restricted form of inclusions, which is close to the database setting considered by Green, Karvounarakis, and Tannen (2007). We also investigate various restrictions of the general setting. Our results are for *combined complexity*, when both the query and the ontology are considered as the input. Proofs are available in our technical report (Bourgau and Ozaki 2018).

## Attributed DL-Lite

Attributed DLs are defined over the usual DL signature with countable sets of *concept names*  $N_C$ , *role names*  $N_R$ , and *individual names*  $N_I$ . We consider an additional set  $N_U$  of *set variables* and a set  $N_V$  of *object variables*. Annotation sets are defined as finite binary relations, understood as sets of attribute-value pairs. Attributes and values refer to domain elements and are syntactically denoted by individual names. To describe annotation sets, we use *specifiers*. The set  $S$  of specifiers contains the following expressions:

- set variables  $X \in N_U$ ;
- *closed specifiers*  $[a_1 : v_1, \dots, a_n : v_n]$ ; and
- *open specifiers*  $[a_1 : v_1, \dots, a_n : v_n]$ ,

where  $a_i \in N_I$  and  $v_i$  is either an individual name in  $N_I$ , an object variable in  $N_V$ , or an expression of the form  $X.a$ , with  $X$  a set variable in  $N_U$  and  $a$  an individual name in  $N_I$ . We use  $X.a$  to refer to the (finite, possibly empty) set of all values of attribute  $a$  in an annotation set  $X$ . A *ground specifier* is a closed or open specifier that only contains individual names. Intuitively, closed specifiers define specific annotation sets whereas open specifiers merely provide lower bounds (Krötzsch et al. 2017).

**Syntax.** A DL-Lite<sub>R</sub><sup>R</sup> *role* (resp. *concept*) *assertion* is an expression  $R(a, b)@S$  (resp.  $A(a)@S$ ), with  $R \in N_R$  (resp.  $A \in N_C$ ),  $a, b \in N_I$ , and  $S \in S$  a ground closed specifier. DL-Lite<sub>R</sub><sup>R</sup> *role* and *concept inclusions* are of the form  $X : S$  ( $P \sqsubseteq Q$ ) and  $X : S$  ( $B \sqsubseteq C$ ) respectively, where  $X \in N_U$ ,  $S \in S$  is a closed or open specifier, and  $P, Q$  and  $B, C$  are respectively role and concept expressions defined by the following syntax, where  $A \in N_C$ ,  $R \in N_R$  and  $S \in S$ :

$$\begin{aligned} P &::= R@S \mid R^-@S, & Q &::= P \mid \neg P, \\ B &::= A@S \mid \exists P, & C &::= B \mid \neg B. \end{aligned}$$

We further require that all variables are *safe*. For set variables, this means that if  $Y \in N_U$  occurs on the right side of an inclusion (or in a specifier  $S$  such that the prefix of the inclusion is  $X : S$  and  $X$  occurs on the right side), then the specifier of the left side expression is  $Y$ . For object variables, if they occur on the right side of an inclusion then they must also occur on the left side or in  $S$  such that  $X : S$  and  $X$  occurs on the left. Note that if object variables occur in  $S$  with  $X : S$  in the prefix and  $X$  on the right side, then  $X$  is the specifier on the left by the safety definition. If the prefix of an inclusion is

$X : S$  and  $X$  does not occur in the role/concept expressions of the inclusion, we may omit  $X : S$ .

A DL-Lite<sub>R</sub><sup>R</sup> *ontology* is a set of DL-Lite<sub>R</sub><sup>R</sup> assertions, role and concept inclusions. Also, we say that a DL-Lite<sub>R</sub><sup>R</sup> ontology is *ground* if it does not contain variables. To simplify notation, we omit the specifier  $[\ ]$  (meaning “any annotation set”) in role or concept expressions. In this sense, any DL-Lite<sub>R</sub> axiom is also a DL-Lite<sub>R</sub><sup>R</sup> axiom. Moreover, we omit prefixes of the form  $X : [\ ]$ , which state that there is no restriction on  $X$ . The *size* of an ontology  $\mathcal{O}$  (or a query, defined later), which we may denote with  $|\mathcal{O}|$ , is the length of the string that represents it.

**Example 1.** Our running example’s ontology  $\mathcal{O}_{ex}$  expresses that those who are married (role spouse) to someone are married (concept Married), annotated with the same sources from which the information has been extracted (attribute src):

$$\exists \text{spouse}@X \sqsubseteq \text{Married}@[\text{src} : X.\text{src}].$$

The assertion states that Zsa Zsa Gabor was married to Jack Ryan and it is annotated with the sources of this information:

$$\text{spouse}(\text{gabor}, \text{ryan})@[\text{src} : s_1, \text{src} : s_2].$$

**Semantics.** An interpretation  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  of an attributed DL consists of a non-empty domain  $\Delta^{\mathcal{I}}$  and a function  $\cdot^{\mathcal{I}}$ . Individual names  $a \in N_I$  are interpreted as elements  $a^{\mathcal{I}} \in \Delta^{\mathcal{I}}$ . To interpret annotation sets, we use the set  $\Phi^{\mathcal{I}} := \{\Sigma \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}} \mid \Sigma \text{ is finite}\}$  of all finite binary relations over  $\Delta^{\mathcal{I}}$ . Each concept name  $A \in N_C$  is interpreted as a set  $A^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Phi^{\mathcal{I}}$  of elements with annotations, and each role name  $R \in N_R$  is interpreted as a set  $R^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}} \times \Phi^{\mathcal{I}}$  of pairs of elements with annotations. Each element (pair of elements) may appear with multiple different annotations.

$\mathcal{I}$  satisfies a concept assertion  $A(a)@[a_1 : v_1, \dots, a_n : v_n]$  if  $(a^{\mathcal{I}}, \{(a_1^{\mathcal{I}}, v_1^{\mathcal{I}}), \dots, (a_n^{\mathcal{I}}, v_n^{\mathcal{I}})\}) \in A^{\mathcal{I}}$ . Role assertions are interpreted analogously. Expressions with free set or object variables are interpreted using variable assignments  $\mathcal{Z}$  mapping object variables  $x \in N_V$  to elements  $\mathcal{Z}(x) \in \Delta^{\mathcal{I}}$  and set variables  $X \in N_U$  to finite binary relations  $\mathcal{Z}(X) \in \Phi^{\mathcal{I}}$ . For convenience, we also extend variable assignments to individual names, setting  $\mathcal{Z}(a) = a^{\mathcal{I}}$  for every  $a \in N_I$ . A specifier  $S \in S$  is interpreted as a set  $S^{\mathcal{I}, \mathcal{Z}} \subseteq \Phi^{\mathcal{I}}$  of matching annotation sets. We set  $X^{\mathcal{I}, \mathcal{Z}} := \{\mathcal{Z}(X)\}$  for variables  $X \in N_U$ . The semantics of closed specifiers is defined as:

- $[a : v]^{\mathcal{I}, \mathcal{Z}} := \{\{(a^{\mathcal{I}}, \mathcal{Z}(v))\}\}$  where  $v \in N_I \cup N_V$ ;
- $[a : X.b]^{\mathcal{I}, \mathcal{Z}} := \{\{(a^{\mathcal{I}}, \delta) \mid (b^{\mathcal{I}}, \delta) \in \mathcal{Z}(X)\}\}$ ;
- $[a_1 : v_1, \dots, a_n : v_n]^{\mathcal{I}, \mathcal{Z}} := \{\bigcup_{i=1}^n F_i \mid F_i \in [a_i : v_i]^{\mathcal{I}, \mathcal{Z}}\}$ .

$S^{\mathcal{I}, \mathcal{Z}}$  therefore is a singleton set for set variables and closed specifiers. For open specifiers, however, we define  $[a_1 : v_1, \dots, a_n : v_n]^{\mathcal{I}, \mathcal{Z}}$  to be the set:

$$\{F \subseteq \Phi^{\mathcal{I}} \mid F \supseteq G \text{ for } \{G\} = [a_1 : v_1, \dots, a_n : v_n]^{\mathcal{I}, \mathcal{Z}}\}.$$

Now given  $A \in N_C$ ,  $R \in N_R$ , and  $S \in S$ , we define:

$$\begin{aligned} (A@S)^{\mathcal{I}, \mathcal{Z}} &:= \{\delta \mid (\delta, F) \in A^{\mathcal{I}} \text{ for some } F \in S^{\mathcal{I}, \mathcal{Z}}\}, \\ (R@S)^{\mathcal{I}, \mathcal{Z}} &:= \{(\delta, \epsilon) \mid (\delta, \epsilon, F) \in R^{\mathcal{I}} \text{ for some } F \in S^{\mathcal{I}, \mathcal{Z}}\}. \end{aligned}$$

Further DL expressions are defined as usual:  $(R^- @ S)^{\mathcal{I}, \mathcal{Z}} = \{(\gamma, \delta) \mid (\delta, \gamma) \in (R @ S)^{\mathcal{I}, \mathcal{Z}}\}$ ,  $\neg P^{\mathcal{I}, \mathcal{Z}} = (\Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}) \setminus P^{\mathcal{I}, \mathcal{Z}}$ ,  $\exists P^{\mathcal{I}, \mathcal{Z}} = \{\delta \mid \text{there is } (\delta, \epsilon) \in P^{\mathcal{I}, \mathcal{Z}}\}$ ,  $\neg C^{\mathcal{I}, \mathcal{Z}} = \Delta^{\mathcal{I}} \setminus C^{\mathcal{I}, \mathcal{Z}}$ .  $\mathcal{I}$  satisfies a concept inclusion  $X : S \sqsubseteq C$  if, for all variable assignments  $\mathcal{Z}$  that satisfy  $\mathcal{Z}(X) \in S^{\mathcal{I}, \mathcal{Z}}$ , we have  $B^{\mathcal{I}, \mathcal{Z}} \subseteq C^{\mathcal{I}, \mathcal{Z}}$ . Satisfaction of role inclusions is defined analogously. An interpretation  $\mathcal{I}$  satisfies an ontology  $\mathcal{O}$ , or is a *model* of  $\mathcal{O}$ , if it satisfies all of its axioms. As usual,  $\models$  denotes the induced logical entailment relation.

**Example 2** (Example 1 cont'd). Consider an interpretation  $\mathcal{I}$  with domain  $\Delta^{\mathcal{I}} = \{\text{gabor, ryan, src, s}_1, \text{s}_2\}$  and such that  $\mathcal{I}$  maps each individual name to itself and

$$\begin{aligned} \text{spouse}^{\mathcal{I}} &= \{(\text{gabor, ryan}, \{(\text{src, s}_1), (\text{src, s}_2)\})\} \\ \text{Married}^{\mathcal{I}} &= \{(\text{gabor}, \{(\text{src, s}_1), (\text{src, s}_2)\})\}. \end{aligned}$$

The interpretation  $\mathcal{I}$  is a model of  $\mathcal{O}_{\text{ex}}$ .

### Reasoning in DL-Lite $_{\text{R}}^{\text{R}}$

In this section, we study the complexity of satisfiability and query answering over DL-Lite $_{\text{R}}^{\text{R}}$  ontologies. Our first result is that the satisfiability problem, which is in NL for DL-Lite $_{\text{R}}$  (Artale et al. 2009), is harder for DL-Lite $_{\text{R}}^{\text{R}}$ . The proof is by reduction from the word problem for polynomially space bounded deterministic Turing Machines (DTM). Annotations raise the complexity because they can encode configurations of a DTM, using expressions of the form  $X.b$  to encode the synchronization of successive configurations.

**Theorem 1.** In DL-Lite $_{\text{R}}^{\text{R}}$ , satisfiability is PSPACE-hard.

To prove the PSPACE upper bound for satisfiability, we use *grounding* (Krötzsch et al. 2017), which is a classical technique that consists in eliminating variables from an ontology to transform it into an equisatisfiable ground ontology. The ground ontology can then be translated into an equisatisfiable DL-Lite $_{\text{R}}$  ontology. The grounding leads to an exponential blowup of the ontology while the translation to DL-Lite $_{\text{R}}$  is polynomial. Since satisfiability of DL-Lite $_{\text{R}}$  ontologies is in NL (Artale et al. 2009), it follows (by (Savitch 1970)) that satisfiability of DL-Lite $_{\text{R}}^{\text{R}}$  ontologies is in PSPACE.

**Theorem 2.** In DL-Lite $_{\text{R}}^{\text{R}}$ , satisfiability is in PSPACE.

We now turn our attention to the problem of querying DL-Lite $_{\text{R}}^{\text{R}}$  ontologies. In the following we only define and deal with conjunctive queries without free variables, i.e., boolean conjunctive queries (BCQ), as the problem of finding certain answers to a query is reducible to BCQ entailment.

**Definition 1** (Attributed Queries). An attributed boolean conjunctive query (BCQ $_{\text{A}}$ )  $q$  is an expression of the form:

$$\exists \mathbf{x}. X_1 : S_1, \dots, X_n : S_n \varphi(\mathbf{x}) \quad (1)$$

where, for  $1 \leq i \leq n$ ,  $X_i$  are the set variables occurring in  $\varphi(\mathbf{x})$ ,  $S_i \in \mathbf{S}$ , and  $\varphi(\mathbf{x})$  is a conjunction of atoms of the form  $A(t) @ S$  or  $R(t, u) @ S$ , with  $A \in \mathbf{N}_C$ ,  $R \in \mathbf{N}_R$ ,  $S \in \mathbf{S}$ , and  $t, u$  individual names in  $\mathbf{N}_I$  or variables in  $\mathbf{x} \subseteq \mathbf{N}_V$ .

We may write  $E(t) @ S$  to refer to an atom of any of the two forms ( $E \in \mathbf{N}_C \cup \mathbf{N}_R$  and  $t$  is a tuple of elements from  $\mathbf{N}_I \cup \mathbf{x}$  of the arity of  $E$ ).

An interpretation  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  satisfies a BCQ $_{\text{A}}$   $q$ , written  $\mathcal{I} \models q$ , if there exists a variable assignment  $\mathcal{Z}$  such that:

- $\mathcal{Z}(X_i) \in S_i^{\mathcal{I}, \mathcal{Z}}$  for all  $1 \leq i \leq n$ ; and
- $(\mathcal{Z}(t), F) \in E^{\mathcal{I}}$  for some  $F \in S^{\mathcal{I}, \mathcal{Z}}$ , for every atom  $E(t) @ S$  occurring in  $q$ .

A BCQ $_{\text{A}}$   $q$  is entailed by  $\mathcal{O}$ , written  $\mathcal{O} \models q$ , iff  $q$  is satisfied by every model of  $\mathcal{O}$ . A BCQ $_{\text{A}}$  that consists of a single atom is an attributed boolean atomic query (BAQ $_{\text{A}}$ ). We say that a BCQ $_{\text{A}}$  is ground if it contains only ground specifiers.

BCQ $_{\text{A}}$  can express conditions on annotations, for instance require that there exists an annotation set where a given attribute is present or has a specific value.

**Example 3.** We modify  $\mathcal{O}_{\text{ex}}$  to express that those who have a spouse are married, associated with the same annotations:

$$\exists \text{spouse} @ X \sqsubseteq \text{Married} @ X.$$

We also add assertions stating that Zsa Zsa Gabor was married to Jack Ryan from 1975 to 1976, while Liz Taylor was married to Richard Burton from 1975 to 1976, as well as the sources of this information:

$$\begin{aligned} &\text{spouse}(\text{gabor, ryan}) @ [\text{start} : 1975, \text{end} : 1976, \text{src} : \text{s}_1], \\ &\text{spouse}(\text{gabor, ryan}) @ [\text{start} : 1975, \text{end} : 1976, \text{src} : \text{s}_2], \\ &\text{spouse}(\text{taylor, burton}) @ [\text{start} : 1975, \text{end} : 1976, \text{src} : \text{s}_3]. \end{aligned}$$

The following query expresses that Gabor and Taylor were married (to someone) with the same start and end dates:

$$\begin{aligned} q_{\text{ex}} &= \exists xy \text{ Married}(\text{gabor}) @ [\text{start} : x, \text{end} : y] \wedge \\ &\quad \text{Married}(\text{taylor}) @ [\text{start} : x, \text{end} : y]. \end{aligned}$$

By the semantics of DL-Lite $_{\text{R}}^{\text{R}}$ , it follows that  $\mathcal{O}_{\text{ex}} \models q_{\text{ex}}$ . This other query expresses that a set of sources that includes  $\text{s}_1$  and is associated with Gabor's married status is also associated with Taylor's married status:

$$\begin{aligned} q'_{\text{ex}} &= X : [\text{src} : \text{s}_1] \text{ Married}(\text{gabor}) @ X \wedge \\ &\quad \text{Married}(\text{taylor}) @ [\text{src} : X.\text{src}]. \end{aligned}$$

By the semantics of DL-Lite $_{\text{R}}^{\text{R}}$ , it follows that  $\mathcal{O}_{\text{ex}} \not\models q'_{\text{ex}}$ .

While BCQ entailment is NP-complete in DL-Lite $_{\text{R}}$ , it follows from Theorem 1 that BAQ $_{\text{A}}$  entailment is already PSPACE-hard. Indeed, satisfiability can be reduced to BAQ $_{\text{A}}$  entailment:  $\mathcal{O}$  is unsatisfiable iff  $\mathcal{O} \models A(a)$  where  $A$  and  $a$  are respectively a concept and an individual name that do not occur in  $\mathcal{O}$ . We show PSPACE-completeness of BCQ $_{\text{A}}$  entailment by describing how to decide  $\mathcal{O} \models q$  for a BCQ $_{\text{A}}$   $q$ , using only polynomial space w.r.t. the size of  $\mathcal{O}$  and  $q$ .

The main ingredients to prove our result are grounding, translation to DL-Lite $_{\text{R}}$ , and also *query rewriting*, a prominent query answering technique for DL-Lite $_{\text{R}}$  in which the query is rewritten w.r.t. the concept and role inclusions, to be evaluated over the assertions as in the classical database setting. However, as the ground version of  $\mathcal{O}$  is of exponential size and the number of rewritten queries is exponential, we do not compute them but instead guess the DL-Lite $_{\text{R}}$  translation  $\text{dl}(q_{\mathcal{Z}})$  of a grounded version  $q_{\mathcal{Z}}$  of  $q$  together with one of its rewritings  $q'$ . We can verify in NP that  $q'$  is entailed by the DL-Lite $_{\text{R}}$  translation of the assertions of  $\mathcal{O}$ , in PTIME that  $\text{dl}(q_{\mathcal{Z}})$  is the DL-Lite $_{\text{R}}$  translation of a grounded version of  $q$ ,

and in PSPACE that  $q'$  is indeed a rewriting of  $\text{dl}(q_Z)$ . For this last step, we propose a non-deterministic adaptation of the rewriting algorithm PerfectRef for DL-Lite $_{\mathcal{R}}$  by Calvanese et al. (2007) that takes as input  $\text{dl}(q_Z)$ ,  $q'$  and  $\mathcal{O}$ . The main idea is to rewrite  $\text{dl}(q_Z)$  by guessing at each step an atom of the query together with a positive inclusion that would appear in the DL-Lite $_{\mathcal{R}}$  translation of the grounding of  $\mathcal{O}$ , thus avoiding the computation of the grounding of  $\mathcal{O}$ .

**Theorem 3.** *In DL-Lite $_{\mathcal{Q}}^{\mathcal{R}}$ , BCQ $_{\mathcal{Q}}$  entailment is in PSPACE.*

The result of Theorem 3, which is for combined complexity, contrasts with the EXPTIME-hardness w.r.t. *data complexity* (only w.r.t. the data size) for MARPL, an attributed logic based on Datalog (Marx, Krötzsch, and Thost 2017). Finally, we show lower complexity bounds in the case of *ground* ontologies. Indeed, when  $\mathcal{O}$  is ground, one can build a DL-Lite $_{\mathcal{R}}$  ontology of polynomial size w.r.t. the size of  $\mathcal{O}$  that entails the DL-Lite $_{\mathcal{R}}$  translation of a grounded version of  $q$  if and only if  $\mathcal{O} \models q$ .

**Theorem 4.** *For ground DL-Lite $_{\mathcal{Q}}^{\mathcal{R}}$  ontologies, satisfiability is in PTIME and BCQ $_{\mathcal{Q}}$  entailment is NP-complete.*

### Querying Using Provenance Semirings

In this section, we investigate attributed DL in light of provenance semirings (Green, Karvounarakis, and Tannen 2007) and enhance the semantics of DL-Lite $_{\mathcal{Q}}^{\mathcal{R}}$  to deal with provenance information. Semirings generalize formalisms such as why-provenance, lineages used in view maintenance, or the lineage used by the Trio uncertain management system (Senellart 2017). The main motivation is to use annotations to answer questions such as “Where does the result come from?”. Assuming that facts are annotated with their sources, we want to know which combinations of sources lead to the entailment of a query. Such annotations may represent various types of information, such as trust, probability, multiplicity or data classification (see Example 8).

**Example 4** (Example 3 cont’d). *The result of the query  $q_{\text{ex}}$  over the ontology  $\mathcal{O}_{\text{ex}}$  can be obtained from source  $s_3$  together with any of  $s_1, s_2$ . Provenance semirings can formalize this information in the form of a provenance polynomial:*

$$(s_1 + s_2) \times s_3.$$

The intuitive meaning is that  $+$  corresponds to *alternative* use of data and  $\times$  to *joint* use of data. The goal of this section is to embed the formalism of provenance semirings into the semantics of DL-Lite $_{\mathcal{Q}}^{\mathcal{R}}$ , so that we can associate annotations using provenance polynomials to queries (e.g., associate the annotation  $\text{src} : (s_1 + s_2) \times s_3$  to the query  $q_{\text{ex}}$  of Example 3).

We define DL-Lite $_{\mathcal{Q}, \mathbb{K}}^{\mathcal{R}}$  as an order-sorted version of DL-Lite $_{\mathcal{Q}}^{\mathcal{R}}$ . Elements of different sorts correspond to sets of *individual names*  $N_I$ , *provenance sums*  $N_S$  and *provenance polynomials*  $N_P$ . We represent provenance polynomials with the *positive algebra provenance semiring* for  $N_I$ , defined as the commutative semiring of polynomials with variables in  $N_I$  and coefficients from  $\mathbb{N}$ , with operations defined as usual:  $\mathcal{K} = (\mathbb{N}[N_I], +, \times, 0, 1)$ . We denote by  $N_P$  the set of polynomials of  $\mathcal{K}$  and by  $N_S$  the subset of  $N_P$  containing the sums of the commutative monoid  $(\mathbb{N}[N_I], +, 0)$ . We thus have

$N_I \subseteq N_S \subseteq N_P$ . We may use the symbols  $\sum$  and  $\prod$  to denote sum and product of elements in  $N_P$  (which will then also be in  $N_P$ ). Elements of  $N_S$  are used as values in the ontology specifiers while elements of  $N_P$  appear as values in the query specifiers. Non-linear polynomials indicate the use of several assertions to derive a query, while provenance sums indicate that a query can be derived from different sources.

Role and concept inclusions in DL-Lite $_{\mathcal{Q}, \mathbb{K}}^{\mathcal{R}}$  are defined similarly as in DL-Lite $_{\mathcal{Q}}^{\mathcal{R}}$ , with the only difference that we allow elements from  $N_S$  to be values of attributes in the specifiers. Concept and role assertions are defined as in DL-Lite $_{\mathcal{Q}}^{\mathcal{R}}$ . The fact that we do not allow values from  $N_S$  in the assertions does not change the expressivity of DL-Lite $_{\mathcal{Q}, \mathbb{K}}^{\mathcal{R}}$ , since inclusions can enforce the entailments of such assertions.

**Example 5.** *The following concept inclusion restricts that of Example 3 by further requiring that the fact that someone has a spouse has to be associated both with  $s_1$  and with  $s_2$  to conclude that this person is married.*

$$X : [\text{src} : s_1 + s_2] \quad (\exists \text{spouse} @ X \sqsubseteq \text{Married} @ X)$$

**Provenance Semantics.** We now introduce the semantics of DL-Lite $_{\mathcal{Q}, \mathbb{K}}^{\mathcal{R}}$ , based on provenance sums. A *provenance-interpretation*  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  is such that  $\cdot^{\mathcal{I}}$  maps polynomials  $a$  and  $b$  in  $N_P$  to the same element  $a^{\mathcal{I}} = b^{\mathcal{I}}$  in  $\Delta^{\mathcal{I}}$  if and only if they are mathematically equal<sup>2</sup>. We denote by  $\Delta_I^{\mathcal{I}}$  the domain of individuals and by  $\Delta_S^{\mathcal{I}}$  the domain of provenance sums, which are the subsets of  $\Delta^{\mathcal{I}}$  corresponding to the image of elements in  $N_I$  and  $N_S$ , respectively. Thus  $\Delta_I^{\mathcal{I}} \subseteq \Delta_S^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$ . To capture the semantics of provenance sums we develop a notion of closure. Intuitively, if a fact is annotated with  $n$  sources then it should also be annotated with the sum of *any subset* of these sources, since the fact can be retrieved alternatively by any source from this subset. For example, assume  $(a, F_1), \dots, (a, F_n)$  are in the interpretation of a concept or a role name  $E$ . If there is  $(\text{src}^{\mathcal{I}}, s_i^{\mathcal{I}})$  in each  $F_i$  and these annotation sets only differ by such pairs, then for each subset of  $\{s_1, \dots, s_n\}$ , the interpretation of  $E$  should have  $(a, F_s)$  with  $F_s$  differing from  $F_i$  only by the pair  $(\text{src}^{\mathcal{I}}, s^{\mathcal{I}})$ , where  $s$  is the sum of the elements of the subset.

We say that  $G, H \in \Phi^{\mathcal{I}}$  are *differentiated* by  $p$  in  $F$  if

$$F = G \setminus \{(p, a) \mid (p, a) \in G\} = H \setminus \{(p, b) \mid (p, b) \in H\}.$$

In this case, we denote by  $G +^p H$  the set

$$F \cup \{(p, (a + b)^{\mathcal{I}}) \mid \{a, b\} \subseteq N_P, (p, a^{\mathcal{I}}) \in G, (p, b^{\mathcal{I}}) \in H\}.$$

A sum of possibly more than two annotation sets differentiated by  $p$  may be denoted by  $\sum_{1 \leq i \leq n} G_i$  and is unique by the commutative law. For  $E \in N_C \cup N_R$ , and  $\mathbf{a}$  a tuple of the arity of  $E$ , we say that  $G +^p H$  is *non-primitive* for  $\mathbf{a}$  and  $E^{\mathcal{I}}$  if  $\{(\mathbf{a}, G), (\mathbf{a}, H)\} \subseteq E^{\mathcal{I}}$ . We denote by  $E_{\mathcal{I}, \mathbf{a}, F}^p$  the set of annotation sets  $G$  pairwise differentiated by  $p$  in  $F \in \Phi^{\mathcal{I}}$  such that  $(\mathbf{a}, G) \in E^{\mathcal{I}}$  with  $G$  primitive for  $\mathbf{a}$  and  $E^{\mathcal{I}}$ .

**Definition 2** (Closure of  $E^{\mathcal{I}}$ ).  *$E^{\mathcal{I}}$  is closed under sum if for all tuples  $\mathbf{a}$  (in  $\Delta^{\mathcal{I}}$  or  $\Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$  according to the arity of  $E$ ),*

<sup>2</sup>According to associative, commutative and distributive laws. E.g.,  $(a + b)^{\mathcal{I}} = (b + a)^{\mathcal{I}}$  by the commutative law.

$\{(\alpha, \sum_{G \in \sigma} G) \mid \sigma \subseteq E_{\mathcal{I}, \alpha, F}^p, \sigma \neq \emptyset\} \subseteq E^{\mathcal{I}}$  for every  $p \in \Delta^{\mathcal{I}}$  and every  $F \in \Phi^{\mathcal{I}}$ .

A provenance-interpretation  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  is *well-founded* if  $E^{\mathcal{I}}$  is closed under sum for all  $E \in N_C \cup N_R$ . For all  $E \in N_C \cup N_R$  and  $\alpha$  with elements in  $\Delta^{\mathcal{I}}$ , we also require that the *support* of  $E^{\mathcal{I}}$  and  $\alpha$  defined as  $\{F \mid (\alpha, F) \in E^{\mathcal{I}}\}$  is finite. This ensures that the sum in Definition 2 is finite. An interpretation of DL-Lite $_{\mathbb{K}}^{\mathcal{R}}$  is a well-founded provenance-interpretation. We denote by  $S_5$  the set of specifiers defined in the same way as  $S$  except that we use  $N_5$  instead of  $N_I$  when defining values of attributes. The semantics of specifiers in  $S_5$  is defined as expected following the definition given in the Section ‘Attributed DL-Lite’ and we use the same notions of satisfiability and entailment. In Definition 2 we consider all subsets of  $E_{\mathcal{I}, \alpha, F}^p$  rather than the sum of its elements. This is to ensure monotonicity of DL-Lite $_{\mathbb{K}}^{\mathcal{R}}$ . Otherwise, given for example  $A(a)@[p : a]$  and  $A(a)@[p : b]$  we would lose the entailment  $A(a)@[p : a + b]$  by adding  $A(a)@[p : c]$ .

**Example 6.** Consider the ontology  $\mathcal{O}$  with the assertions  $\text{spouse}(\text{gabor}, \text{ryan})@[\text{src} : s_1]$ ,  $\text{spouse}(\text{gabor}, \text{ryan})@[\text{src} : s_2]$  and the concept inclusion of Example 5. Let  $\mathcal{I}$  have domain  $\Delta^{\mathcal{I}} = \{\text{gabor}, \text{ryan}, \text{src}, s_1, s_2, s_1 + s_2\}$ , interpret each individual name by itself,  $(s_1 + s_2)^{\mathcal{I}} = s_1 + s_2$ , and

$$\text{spouse}^{\mathcal{I}} = \{(\text{gabor}, \text{ryan}, G), (\text{gabor}, \text{ryan}, H), (\text{gabor}, \text{ryan}, G +^{\text{src}} H)\}$$

$$\text{Married}^{\mathcal{I}} = \{(\text{gabor}, G +^{\text{src}} H)\} \text{ where } G = \{(\text{src}, s_1)\}, H = \{(\text{src}, s_2)\} \text{ and } G +^{\text{src}} H = \{(\text{src}, s_1 + s_2)\}.$$

$\text{spouse}^{\mathcal{I}}$  and  $\text{Married}^{\mathcal{I}}$  are closed under sum,  $\mathcal{I}$  is a model of  $\mathcal{O}$  and  $\mathcal{O} \models \text{spouse}(\text{gabor}, \text{ryan})@[\text{src} : s_1 + s_2]$ .

We denote by  $S_P$  the set of specifiers defined in the same way as  $S_5$  except that we use  $N_P$  instead of  $N_5$  for the values of attributes. The semantics of specifiers in  $S_P$  is as expected from the Section ‘Attributed DL-Lite’. We assume that all polynomials occurring in a specifier in  $S_P$  are of the form  $\sum_{1 \leq i \leq n_1} \prod_{1 \leq j \leq n_2} a_{i,j}$ , where all  $a_{i,j} \in N_I$ . Given an interpretation  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  and  $\{F, G\} \subseteq \Phi^{\mathcal{I}}$ , let  $F \times G$  be:

$$\{(p, (a \times b)^{\mathcal{I}}) \mid \{a, b\} \subseteq N_P, (p, a^{\mathcal{I}}) \in F, (p, b^{\mathcal{I}}) \in G\}.$$

Unlike  $+^p$ ,  $\times$  is not parameterized by an attribute because products combine different information, whereas sums represent alternative ways of obtaining the same information (i.e., tuple plus the same other attribute-value pairs). A product of annotation sets may be denoted by  $\prod_{1 \leq i \leq n} G_i$ . We next define semiring attributed queries, which allow a ground specifier to be associated to the whole conjunction of atoms.

**Definition 3** (Semiring Attributed Queries). A semiring attributed boolean conjunctive query ( $BCQ_{\mathbb{K}}^{\mathcal{R}}$ ) is an expression of the form:

$$\exists x. X_1 : S_1, \dots, X_n : S_n (\varphi(x))@S,$$

where  $S$  is a ground specifier in  $S_P$ , for  $1 \leq i \leq n$ ,  $X_i \in N_U$  are the set variables occurring in  $\varphi(x)$  and  $S_i \in S_5$ , and

$$\varphi(x) = \bigwedge_{1 \leq j \leq m} E_j(t_j)@T_j$$

where for  $1 \leq j \leq m$ ,  $T_j \in S_5$ ,  $E_j \in N_C \cup N_R$  and  $t_j$  is a tuple of elements from  $N_I \cup N_R$ .

If  $S = \perp$ , we say that the  $BCQ_{\mathbb{K}}^{\mathcal{R}}$  is plain.

Given a  $BCQ_{\mathbb{K}}^{\mathcal{R}}$   $q$ , let  $q'$  be the  $BCQ_{\mathbb{K}}^{\mathcal{R}}$  that results from removing the outer specifier from  $q$ . Let  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  be an interpretation and let  $\nu_{\mathcal{I}}(q')$  be the set of all variable assignments  $\mathcal{Z}$  that fulfill the conditions of Definition 1 for  $\mathcal{I} \models q'$ .  $\mathcal{I}$  satisfies  $q$ , written  $\mathcal{I} \models q$ , if there is a non-empty  $\chi \subseteq \nu_{\mathcal{I}}(q')$  such that:

1. for any  $\mathcal{Z}, \mathcal{Z}' \in \chi$ , there exists  $X \in N_U$  occurring in  $q$  such that  $\mathcal{Z}(X) \neq \mathcal{Z}'(X)$  or there exists  $x \in x$  such that  $\mathcal{Z}(x) \neq \mathcal{Z}'(x)$ ;
2. for each  $\mathcal{Z} \in \chi$  and  $1 \leq j \leq m$ , we have that  $(\mathcal{Z}(t_j), F_j^{\mathcal{Z}}) \in E_j^{\mathcal{I}, \mathcal{Z}}$  for some  $F_j^{\mathcal{Z}} \in T_j^{\mathcal{I}, \mathcal{Z}}$ ;
3. there is  $p \in \Delta^{\mathcal{I}}$  and  $G \in \Phi^{\mathcal{I}}$  such that all  $H^{\mathcal{Z}} = \prod_{1 \leq j \leq m} F_j^{\mathcal{Z}}$  with  $\mathcal{Z} \in \chi$  are differentiated by  $p$  in  $G$ , and  $\sum_{\mathcal{Z} \in \chi} H^{\mathcal{Z}} \in S^{\mathcal{I}, \mathcal{Z}}$ .

Essentially, Definition 3 says that: (1) there are different variable assignments which (2) satisfy the homomorphic conditions and (3) correspond to the interpretation of the outer specifier. Our semiring attributed queries can be easily extended so that the outer specifier has fresh and free object variables. In this case the answer to the query would be the set of provenance polynomials related with the respective attribute and the query. Semiring attributed queries can be used to query a DL-Lite $_{\mathbb{K}}^{\mathcal{R}}$  ontology using provenance polynomials, as we illustrate with the following example.

**Example 7** (Example 3 cont’d). We now modify  $q_{\text{ex}}$ , so that we impose provenance constraints on the result:

$$\begin{aligned} &\exists xy (\text{Married}(\text{gabor})@[\text{start} : x, \text{end} : y] \wedge \\ &\quad \text{Married}(\text{taylor})@[\text{start} : x, \text{end} : y]@[\text{src} : \gamma] \\ &\quad \text{where } \gamma \text{ is the polynomial } (s_1 \times s_3) + (s_2 \times s_3) \end{aligned}$$

By the semantics of DL-Lite $_{\mathbb{K}}^{\mathcal{R}}$ , it follows that  $\mathcal{O}_{\text{ex}} \models q_{\text{ex}}$ .

All shared attributes are taken into account when combining the annotations, while the non-shared attributes are irrelevant and lost in the product.

**Example 8.** The query  $(\text{Married}(a) \wedge \text{Married}(b))@S$  with  $S = [\text{src} : s_1 \times s_2, \text{classif} : \text{public} \times \text{confid}, \text{mult} : 2 \times 3]$  is entailed by  $\{\text{Married}(a)@[\text{src} : s_1, \text{classif} : \text{public}, \text{mult} : 2], \text{spouse}(b, c)@[\text{src} : s_2, \text{classif} : \text{confid}, \text{mult} : 3, \text{time} : t]\}$  and the inclusion of Example 3.

The fact that  $a$  and  $b$  are both married is obtained by combining sources  $s_1$  and  $s_2$ , and by having access to both public and confidential information. Note that using inclusions to propagate annotations allows the query derived from assertions with multiplicities 2 and 3 to have multiplicity  $2 \times 3$ , as it would be under the bag semantics (Nikolaou et al. 2017).

When interpreted over provenance-interpretations, ontologies in the DL-Lite $_{\mathbb{K}}^{\mathcal{R}}$  fragment of DL-Lite $_{\mathbb{K}}^{\mathcal{R}}$  (i.e., without sums) can entail queries with sums, as in Example 9.

**Example 9.** Let  $\mathcal{O}$  be the DL-Lite $_{\mathbb{K}}^{\mathcal{R}}$  ontology

$$\{A(a)@[p : a], A(a)@[p : b], A@X \sqsubseteq ER@X\}.$$

Then the query  $\exists xy (R(x, y)@[p : a + b])@[]$  follows from  $\mathcal{O}$  only under the semirings-based semantics.

## Reasoning in DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$

Unfortunately, Theorem 5 shows that provenance sums increase the complexity of the satisfiability problem. The proof is by reduction from the word problem for a polynomially space bounded Alternating Turing Machine (ATM) which is EXPTIME-hard (Chandra, Kozen, and Stockmeyer 1981).

**Theorem 5.** *In DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$ , satisfiability is EXPTIME-hard.*

The hardness result of Theorem 5 holds even for DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontologies without expressions of the form  $\exists P$ , where  $P$  is a role expression. Motivated by this negative result, we investigate restricted cases for query answering. We first show that for the class of DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontologies which do not contain inclusions with expressions of the form  $\exists P$  on the right side, we can check the entailment of  $\text{BCQ}_{\mathcal{O}, \mathbb{K}}$  via a transformation to ground and plain  $\text{BCQ}_{\mathcal{O}, \mathbb{K}}$ s. Given such a DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontology  $\mathcal{O}$ , one can translate a  $\text{BCQ}_{\mathcal{O}, \mathbb{K}}$   $q$  into a set of ground and plain  $\text{BCQ}_{\mathcal{O}, \mathbb{K}}$ s  $\text{gr\_plain}(\mathcal{O}, q)$  such that  $\mathcal{O} \models q$  iff there is some  $q_{\text{gp}} \in \text{gr\_plain}(\mathcal{O}, q)$  that is entailed by an equisatisfiable ground ontology.

We can assume w.l.o.g. that if  $E_j(t_j)@T_j$  occurs in  $q$  then  $T_j \in \mathcal{N}_U$ : if  $T_j$  is a specifier one can always replace it by a fresh  $X \in \mathcal{N}_U$  and add  $X : T_j$  to the prefix of  $q$ , that is:

$$q = \exists x. X_1 : S_1, \dots, X_m : S_m \left( \bigwedge_{1 \leq j \leq m} E_j(t_j)@X_j \right) @ S.$$

Assume  $\star \in \mathcal{N}_I$  does not occur in  $\mathcal{O}$  nor in  $q$  and let  $\mathcal{N}_{\text{Pmin}}$  be a fixed but arbitrary minimal subset of  $\mathcal{N}_{\text{P}}$  such that for each  $a \in \mathcal{N}_{\text{P}}$ ,  $\mathcal{N}_{\text{Pmin}}$  contains an element  $b$  such that  $a$  is mathematically equal to  $b$ . Let  $\mathcal{I}$  be a DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  interpretation with domain  $\Delta^{\mathcal{I}} = \mathcal{N}_{\text{Pmin}}$  and such that  $a^{\mathcal{I}} = a$  for every  $a \in \mathcal{N}_{\text{Pmin}}$ . We say that a variable assignment  $\mathcal{Z}$  is *compatible with  $q$*  if  $\mathcal{Z}(X_j) \in S_j^{\mathcal{I}, \mathcal{Z}}$ ,  $1 \leq j \leq m$ . Let  $q'$  be the result of removing the outer specifier from  $q$ . Given a compatible  $\mathcal{Z}$ , a  $\mathcal{Z}$ -image  $\bigwedge_{1 \leq j \leq m} E_j(t_j)@T_j$  of  $q'$  is obtained by:

- replacing each  $X_j$  with  $T_j = [a : b \mid (a, b) \in \mathcal{Z}(X_j)]$ ;
- replacing each object variable  $x$  by  $\mathcal{Z}(x)$ ;
- if  $\star$  occurs in some  $T_j$ , replacing  $\star$  by  $\star_{\mathcal{T}_j}$ , where  $\mathcal{T}_j$  is the set of attribute-value pairs in  $T_j$  that do not contain  $\star$ .

Given a ground specifier  $T$ , let  $F_T := \{(a^{\mathcal{I}}, b^{\mathcal{I}}) \mid a : b \text{ occurs in } T\} \in \Phi^{\mathcal{I}}$ . We define  $\text{gr\_plain}(\mathcal{O}, q)$  as the set of ground plain  $\text{BCQ}_{\mathcal{O}, \mathbb{K}}$ s:

$$q_{\text{gp}} = \left( \bigwedge_{1 \leq i \leq n} \left( \bigwedge_{1 \leq j \leq m} E_j(t_j^i)@S_j^i \right) \right) @ [ ]$$

where the annotation sets  $F_i = \prod_{1 \leq j \leq m} F_{S_j^i}$  ( $1 \leq i \leq n$ ) are such that there exists  $p$  such that (i) the  $F_i$  are differentiated by  $p$  in some annotation set, (ii) each  $F_i$  contains some  $(p, a)$  with  $a \in \mathcal{N}_{\text{P}}$ , and (iii)  $\sum_{1 \leq i \leq n} F_i \in S^{\mathcal{I}}$ . Also,  $\bigwedge_{1 \leq j \leq m} E_j(t_j^i)@S_j^i$  is a  $\mathcal{Z}$ -image of  $q'$  with attribute-value pairs built from elements of  $\mathcal{N}_S$ . By construction,  $q_{\text{gp}}$  does not contain variables.

**Example 10** (Example 3 cont'd). *The query below is a ground and plain version of the query in Example 7 which is*

*entailed by  $\mathcal{O}_{\text{ex}}$ .*

Married(gabor)@[start : 1975, end : 1976, src :  $s_1 + s_2$ ]  $\wedge$   
Married(taylor)@[start : 1975, end : 1976, src :  $s_3$ ].

One can show that, for DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontologies  $\mathcal{O}$  without expressions of the form  $\exists P$  on the right side of inclusions,  $\mathcal{O} \models q$  iff there is  $q_{\text{gp}} \in \text{gr\_plain}(\mathcal{O}, q)$  such that  $\mathcal{O}_{\text{gr}} \models q_{\text{gp}}$ , where  $\mathcal{O}_{\text{gr}}$  is an equisatisfiable ground ontology, obtained in a way similar to our construction of  $\text{gr\_plain}(\mathcal{O}, q)$  but imposing that the image of the variable assignments is over a finite set of individual names defined in terms of  $\mathcal{O}$ . In the case where  $\mathcal{O}$  is ground, we further have a polynomial bound on the size of such  $q_{\text{gp}}$ .

**Lemma 1.** *Let  $q$  be a  $\text{BCQ}_{\mathcal{O}, \mathbb{K}}$  and let  $\mathcal{O}$  be a ground DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontology without expressions of the form  $\exists P$  on the right side of inclusions.  $\mathcal{O} \models q$  iff there is  $q_{\text{gp}} \in \text{gr\_plain}(\mathcal{O}, q)$  such that (i)  $\mathcal{O}_{\text{gr}} \models q_{\text{gp}}$ , (ii) the size of  $q_{\text{gp}}$  is polynomial in  $|q|$  and  $|\mathcal{O}|$  and (iii) deciding  $q_{\text{gp}} \in \text{gr\_plain}(\mathcal{O}, q)$  is in PTIME.*

Lemma 1 does not hold for arbitrary DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontologies, as illustrated by Example 11.

**Example 11.** *Let  $\mathcal{O}$  be the DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontology  $\{A \sqsubseteq \exists R, \exists R^- \sqsubseteq A@[p : b], \exists R^- \sqsubseteq \neg B, B(a), A(a)@[p : b]\}$ . Then,  $\mathcal{O}$  entails  $q = \exists x(A(x)@[p : b + b])$ , since there would be an  $R$ -successor in the anonymous part of the model, but there is no  $q_{\text{gp}} \in \text{gr\_plain}(\mathcal{O}, q)$  such that  $\mathcal{O} \models q_{\text{gp}}$ .*

We now use the polynomial bound in Lemma 1 to show an upper bound for a fragment, called *simple*, where we only allow inclusions of the form  $E_1@S \sqsubseteq E_2@T$ , with  $E_1$  and  $E_2$  concept/role names and  $S$  and  $T$  ground specifiers. We establish the complexity of  $\text{BCQ}_{\mathcal{O}, \mathbb{K}}$  entailment from simple ontologies. This case is close to the classical problem of query answering over databases, considered by Green, Karvounarakis, and Tannen (2007). Theorem 6 states that this complexity remains the same as in the database case.

**Theorem 6.**  *$\text{BCQ}_{\mathcal{O}, \mathbb{K}}$  entailment from a simple DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontology is NP-complete.*

*Proof.* Let  $\mathcal{O}$  be a simple DL-Lite $_{\mathcal{O}, \mathbb{K}}^{\mathcal{R}}$  ontology. We first show that one can decide in NP whether  $E(a)@S$  is entailed from  $\mathcal{O}$ , where  $S$  is a ground specifier.

**Claim 1.** Deciding whether  $\mathcal{O} \models E(a)@S$  is in NP.

*Proof of Claim 1* We first guess the set  $\mathcal{Q}$  of all atomic queries of the form  $E_0(a)@T_0$  entailed by  $\mathcal{O}$  such that  $E_0@T_0$  occurs in  $\mathcal{O}$  and an ordering for the entailment of such queries. If  $T_0$  is an open specifier then replace it in  $\mathcal{Q}$  by  $T_{0, \star}$ , defined as the ground closed specifier containing all attribute-value pairs in  $T_0$  plus  $\star_S : \star_S$  with  $S$  the set of attribute-value pairs in  $T_0$ . We make the usual assumption that individual names of the form  $\star_S$  do not occur in  $\mathcal{O}$  and  $E(a)@S$ . Denote by  $\mathcal{Q}_q$  the subset of  $\mathcal{Q}$  containing all atomic queries which precede  $q$  in the ordering. For each guessed query  $q = E_0(a)@T_0$ :

- Denote by  $F_T$  the set  $\{(a, b) \mid a : b \text{ occurs in } T\}$  for any ground specifier  $T$  and let  $E_0(a)@S_1, \dots, E_0(a)@S_n$  be the assertions and atomic queries in  $\mathcal{O} \cup \mathcal{Q}_q$  where  $E_0$  and  $a$  occur.

- Guess a tree of annotation sets rooted either in  $F_{T_0}$  if  $T_0$  is a closed specifier, or in a superset  $F$  of  $F_{T_0}$  if  $T_0$  is an open specifier, where each non-leaf node  $F$  is the parent of children  $G_1, \dots, G_m$  such that  $F = \sum_{1 \leq i \leq m} G_i$ , for some attribute  $p$ , and such that each leaf is either: one of  $F_{S_1}, \dots, F_{S_n}$ , or some  $F_T$  (or  $F_{T_\star}$  if  $T$  is open) such that there exist  $E_1 @ T_1 \sqsubseteq E_0 @ T$  and  $E_1(a) @ T_1$  (or  $E_1(a) @ T_{1,\star}$  if  $T_1$  is open) in  $\mathcal{O} \cup \mathcal{Q}_q$ .

Check in polynomial time whether the trees satisfy the described conditions. The size of  $\mathcal{Q}$  (and so the number of trees to guess and the size of the ordering) is bounded by the number of atomic queries  $E_0(a) @ T_0$  that can be built from concept/role expressions and individual names in  $\mathcal{O}$ , so it is polynomial in the size of  $\mathcal{O}$ .

To check whether  $\mathcal{O} \models E(a) @ S$ , we check whether  $E(a) @ S \in \mathcal{Q}$  (assuming w.l.o.g. that  $E @ S$  occurs in  $\mathcal{O}$ ). The size of each guessed tree is polynomial in the size of  $\mathcal{O}$  since each leaf corresponds to an assertion/atomic query in  $\mathcal{O}$  or  $\mathcal{Q}$  (or an assertion/atomic query in  $\mathcal{O}$  or  $\mathcal{Q}$  together with an inclusion in  $\mathcal{O}$ ) and they do not repeat in the tree. Thus, one can decide whether  $\mathcal{O} \models E(a) @ S$  in NP.

By Lemma 1,  $\mathcal{O} \models q$  iff there exists  $q_{\text{gp}} \in \text{gr\_plain}(\mathcal{O}, q)$  such that  $\mathcal{O}_{\text{gr}} \models q_{\text{gp}}$ . Moreover the size of  $q_{\text{gp}}$  is polynomial in the size of  $q$  and  $\mathcal{O}$  and  $q_{\text{gp}}$  does not contain variables. We thus get the NP upper bound by guessing  $q_{\text{gp}}$  as well as certificates that  $\mathcal{O}_{\text{gr}} \models E(a) @ S$  for each  $E(a) @ S$  in  $q_{\text{gp}}$ , using Claim 1 (indeed,  $\mathcal{O}_{\text{gr}}$  is also a simple ontology and is of polynomial size w.r.t.  $\mathcal{O}$ ). The lower bound comes from the complexity of BCQ entailment in relational databases.  $\square$

One of the difficulties in showing Theorem 6 for arbitrary  $\text{DL-Lite}_{\mathcal{R}, \mathbb{K}}^{\mathcal{R}}$  ontologies is that one can express that elements in the anonymous part of the model are distinct, as illustrated in Example 11, and then our translation does not hold. In this case,  $\text{gr\_plain}(\mathcal{O}, q)$  needs to include queries with inequalities to distinguish anonymous elements, and entailment of BCQs with inequalities over  $\text{DL-Lite}_{\mathcal{R}}^{\mathcal{R}}$  ontologies easily leads to undecidability (e.g., see Theorem 13 in (Gutiérrez-Basulto et al. 2015)).

We now show an upper bound for *satisfiability* in  $\text{DL-Lite}_{\mathcal{R}, \mathbb{K}}^{\mathcal{R}}$  by translating the ontology into an equisatisfiable ontology in a DL that we call  $\text{DL-Lite}_{\text{Horn}}^{\mathcal{R}, \sqcap}$ , which extends  $\text{DL-Lite}_{\mathcal{R}}^{\mathcal{R}}$  with *conjunctions on the left side of concept and role inclusions*. Our translation is double-exponential since in  $\text{DL-Lite}_{\mathcal{R}, \mathbb{K}}^{\mathcal{R}}$  we need to ensure, e.g., that elements in the extension of  $E @ [\text{src}: s_1]$  and  $E @ [\text{src}: s_2]$  are also in the extension of  $E @ [\text{src}: s_1 + s_2]$ .

**Theorem 7.** *In  $\text{DL-Lite}_{\mathcal{R}, \mathbb{K}}^{\mathcal{R}}$ , satisfiability is in 2EXPTIME.*

*Sketch.* We first ground the ontology and then translate it into  $\text{DL-Lite}_{\text{Horn}}^{\mathcal{R}, \sqcap}$ . We encode the semantics of provenance sums using a double-exponential number of concept and role inclusions with conjunctions on the left side. Since satisfiability in  $\text{DL-Lite}_{\text{Horn}}^{\mathcal{R}, \sqcap}$  is in PTIME (Artale et al. 2015) (Theorem 14), the 2EXPTIME upper bound follows.  $\square$

We next analyse entailment of *plain*  $\text{BCQ}_{\mathcal{R}, \mathbb{K}}$  w.r.t.  $\text{DL-Lite}_{\mathcal{R}, \mathbb{K}}^{\mathcal{R}}$  ontologies: the outer specifier is of the form  $\lfloor \rfloor$  but

inner specifiers can contain provenance sums (as in Ex. 9). We use the fact that BCQ entailment in  $\text{DL-Lite}_{\text{Horn}}^{\mathcal{R}, \sqcap}$  is in NP (Calì, Gottlob, and Pieris 2012, proof of Theorem 3.3).

**Theorem 8.** *In  $\text{DL-Lite}_{\text{Horn}}^{\mathcal{R}, \sqcap}$ , BCQ entailment is in NP.*

Theorem 9 establishes an upper bound for plain queries.

**Theorem 9.** *In  $\text{DL-Lite}_{\mathcal{R}, \mathbb{K}}^{\mathcal{R}}$ , entailment of plain  $\text{BCQ}_{\mathcal{R}, \mathbb{K}}$  is in N2EXPTIME.*

*Sketch.* The proof uses the translation to  $\text{DL-Lite}_{\text{Horn}}^{\mathcal{R}, \sqcap}$  which leads to a double-exponential blowup of the ontology. Here, since queries are plain the translation is as for  $\text{BCQ}_{\mathcal{R}, \mathbb{K}}$ . The result then follows from Theorem 8.  $\square$

## Conclusion

We investigated the complexity of satisfiability and query answering in attributed  $\text{DL-Lite}_{\mathcal{R}}$ , for both the semantics introduced by Krötzsch et al. (2017) and a new semantics based on provenance semirings, which allows to embed provenance polynomials into the query. In particular, we show that these problems are PSPACE-complete for the classical semantics and that in the case of simple ontologies, even query answering under the semirings-based semantics has the same complexity as query answering in  $\text{DL-Lite}_{\mathcal{R}}$ . However, satisfiability of general  $\text{DL-Lite}_{\mathcal{R}, \mathbb{K}}^{\mathcal{R}}$  ontologies is EXPTIME-hard.

**Related Work.** Our attributed ontology language differs from  $\text{DL-Lite}_{\mathcal{A}}$  (Calvanese et al. 2006), which allows to associate values to individuals or pairs of individuals, rather than to assertions, through binary or ternary relations called attribute concepts or attribute roles. In particular, while we can use the same attribute name to annotate different assertions about the same individual or pair of individuals, it would be ambiguous in  $\text{DL-Lite}_{\mathcal{A}}$ . For instance, we can express that Liz Taylor was married to Richard Burton from 1964 to 1974 and from 1975 to 1976 with  $\text{spouse}(\text{taylor}, \text{burton}) @ [\text{start} : 1964, \text{end} : 1974]$ ,  $\text{spouse}(\text{taylor}, \text{burton}) @ [\text{start} : 1975, \text{end} : 1976]$ , while in  $\text{DL-Lite}_{\mathcal{A}}$  we would need reification. The query  $\text{spouse}(\text{taylor}, \text{burton}) @ [\text{start} : x, \text{end} : y]$  that returns the start and end dates of the marriages would be more complex (namely, e.g.,  $\exists z \text{spouse}_1(z, \text{taylor}) \wedge \text{spouse}_2(z, \text{burton}) \wedge \text{start}(z, x) \wedge \text{end}(z, y)$ ). Another difference is the use in  $\text{DL-Lite}_{\mathcal{A}}$  of two distinct alphabets and interpretation domains for the individuals and the values, following the distinction made in OWL between objects and values.

Regarding provenance, the topic has been extensively studied for relational databases (Cheney, Chiticariu, and Tan 2009), but has also drawn attention in other settings, e.g., for Datalog (Deutch et al. 2014), Datalog<sup>+/-</sup> (Lukasiewicz et al. 2014), and Semantic Web data, with numerous works proposing provenance models based on semirings for the evaluation of SPARQL queries over annotated RDF, see e.g., (Theoharis et al. 2011; Zimmermann et al. 2012; Geerts et al. 2016). In particular, Zimmermann et al. consider the possibility of having several annotations with different domains (fuzzy, temporal and provenance) and introduce an annotated version of SPARQL that manipulates explicitly



annotations, while most work on provenance only implicitly propagates provenance annotations.

**Future Work.** Our next step will be the study of the data complexity and the design of practical algorithms for querying attributed DL-Lite ontologies. In particular, we would like to extend the classical DL-Lite rewriting approach to the attributed setting to avoid grounding the ontology. For instance, if an ontology only contains inclusions of the form  $E@X \sqsubseteq F@X$ , the rewriting algorithm for DL-Lite<sub>R</sub> could be adapted to rewrite an attributed query where annotations sets are propagated in the rewriting process (e.g.,  $\text{Married}(\text{gabor})@[\text{start} : 1975]$  in Example 3 could be rewritten into  $\exists y \text{ spouse}(\text{gabor}, y)@[\text{start} : 1975]$ ).

**Acknowledgements.** We thank Stefan Borgwardt and an anonymous reviewer for pointing us to useful references. Partially supported by ANR-16-CE23-0007-01 (“DICOS”).

## References

- Artale, A.; Calvanese, D.; Kontchakov, R.; and Zakharyashev, M. 2009. The DL-Lite family and relations. *J. Artif. Intell. Res.* 36:1–69.
- Artale, A.; Kontchakov, R.; Ryzhikov, V.; and Zakharyashev, M. 2015. Tractable interval temporal propositional and description logics. In *Proceedings of AAAI*.
- Baader, F.; Calvanese, D.; McGuinness, D.; Nardi, D.; and Patel-Schneider, P., eds. 2007. *The Description Logic Handbook: Theory, Implementation, and Applications*. Cambridge University Press, second edition.
- Bienvenu, M., and Ortiz, M. 2015. Ontology-mediated query answering with data-tractable description logics. In *Reasoning Web, Tutorial Lectures*, 218–307.
- Bienvenu, M.; Deutch, D.; and Suchanek, F. M. 2012. Provenance for Web 2.0 data. In *Proceedings of Secure Data Management: 9th VLDB Workshop, SDM 2012*.
- Bourgau, C., and Ozaki, A. 2018. Querying attributed DL-Lite ontologies using provenance semirings. Technical Report 02, Free University of Bozen-Bolzano. Available at <https://www.inf.unibz.it/krdp/KRDB%20files/tech-reports/KRDB18-02.pdf>.
- Calì, A.; Gottlob, G.; and Pieris, A. 2012. Towards more expressive ontology languages: The query answering problem. *Artif. Intell.* 193:87–128.
- Calvanese, D.; De Giacomo, G.; Lembo, D.; Lenzerini, M.; Poggi, A.; and Rosati, R. 2006. Linking data to ontologies: The description logic DL-Lite<sub>A</sub>. In *Proceedings of the OWLED\*06 Workshop on OWL: Experiences and Directions*.
- Calvanese, D.; Giacomo, G. D.; Lembo, D.; Lenzerini, M.; and Rosati, R. 2007. Tractable reasoning and efficient query answering in description logics: The DL-Lite family. *J. of Automated Reasoning* 39(3):385–429.
- Chandra, A. K.; Kozen, D. C.; and Stockmeyer, L. J. 1981. Alternation. *J. of the ACM* 28(1):114–133.
- Cheney, J.; Chiticariu, L.; and Tan, W. C. 2009. Provenance in databases: Why, how, and where. *Foundations and Trends in Databases* 1(4):379–474.
- Deutch, D.; Milo, T.; Roy, S.; and Tannen, V. 2014. Circuits for datalog provenance. In *Proceedings of ICDT*.
- Geerts, F.; Unger, T.; Karvounarakis, G.; Fundulaki, I.; and Christophides, V. 2016. Algebraic structures for capturing the provenance of SPARQL queries. *J. ACM* 63(1):7:1–7:63.
- Green, T. J.; Karvounarakis, G.; and Tannen, V. 2007. Provenance semirings. In *Proceedings of PODS*.
- Gutiérrez-Basulto, V.; Ibáñez García, Y.; Kontchakov, R.; and Kostylev, E. V. 2015. Queries with negation and inequalities over lightweight ontologies. *Web Semant.* 35(P4):184–202.
- Hoffart, J.; Suchanek, F. M.; Berberich, K.; and Weikum, G. 2013. YAGO2: A spatially and temporally enhanced knowledge base from wikipedia. *Artif. Intell.* 194:28–61.
- Krötzsch, M.; Marx, M.; Ozaki, A.; and Thost, V. 2017. Attributed description logics: Ontologies for knowledge graphs. In *Proceedings of ISWC*.
- Krötzsch, M.; Marx, M.; Ozaki, A.; and Thost, V. 2018. Attributed description logics: Reasoning on knowledge graphs. In *Proceedings of IJCAI*.
- Lukasiewicz, T.; Martinez, M. V.; Predoiu, L.; and Simari, G. I. 2014. Information integration with provenance on the semantic web via probabilistic datalog+/- . In *URSW 2011-2013, Revised Selected Papers*.
- Marx, M.; Krötzsch, M.; and Thost, V. 2017. Logic on MARS: Ontologies for generalised property graphs. In *Proceedings of IJCAI*.
- Motik, B.; Cuenca Grau, B.; Horrocks, I.; Wu, Z.; Fokoue, A.; and Lutz, C., eds. 2009. *OWL 2 Web Ontology Language: Profiles*. W3C Recommendation. Available at <http://www.w3.org/TR/owl2-profiles/>.
- Nikolaou, C.; Kostylev, E. V.; Konstantinidis, G.; Kaminski, M.; Grau, B. C.; and Horrocks, I. 2017. The bag semantics of ontology-based data access. In *Proceedings of IJCAI*.
- Rodriguez, M. A., and Neubauer, P. 2010. Constructions from dots and lines. *Bulletin of the American Society for Information Science and Technology* 36(6):35–41.
- Savitch, W. J. 1970. Relationships between nondeterministic and deterministic tape complexities. *Journal of Computer and System Sciences* 4(2):177 – 192.
- Senellart, P. 2017. Provenance and probabilities in relational databases. *SIGMOD Record* 46(4):5–15.
- Suciu, D.; Olteanu, D.; Ré, C.; and Koch, C. 2011. *Probabilistic Databases*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers.
- Theoharis, Y.; Fundulaki, I.; Karvounarakis, G.; and Christophides, V. 2011. On provenance of queries on semantic web data. *IEEE Internet Computing* 15(1):31–39.
- Vrandečić, D., and Krötzsch, M. 2014. Wikidata: A free collaborative knowledgebase. *Commun. ACM* 57(10).
- Zimmermann, A.; Lopes, N.; Polleres, A.; and Straccia, U. 2012. A general framework for representing, reasoning and querying with annotated semantic web data. *J. Web Sem.* 11:72–95.